



Computational Biology

Lecture #4: Introductory Biology & Genomics

Bud Mishra
Professor of Computer Science, Mathematics, & Cell Biology
Oct 03 2005

10/18/2005

© Bud Mishra, 2005

L4-1



Areas of Interest to Bioinformatics

10/18/2005

© Bud Mishra, 2005

L4-2



Active Areas of Research

- ◇ **Human Genome Project:**
 - Read 3 billion base pairs in 46 human chromosomes
- ◇ **Polymorphisms and Haplotyping**
 - **SNPs** (Single Nucleotide Polymorphisms): Catalog the single base pair variations occurring about 1 in 800 base pairs of human genome over the entire populations
 - **CNPs** (Copy Number Polymorphisms)
 - **RFLP-Maps** Restriction Fragment Length Polymorphisms

10/18/2005

© Bud Mishra, 2005

L4-3



Active Areas of Research

- ◇ **Transcription Maps:**
 - Identify all (about 30,000 (??)) the genes in the human genome.
 - Particularly interesting are the ones involved in diseases... Cancer may involve about 100 oncogenes and 1000 tumor suppressor genes
- ◇ **Linkage Analysis:**
 - Relate genes (or polymorphic markers) to phenotypes (gene-expression patterns & externally observable traits) by analyzing genomes of a family (kinship) or over a population.

10/18/2005

© Bud Mishra, 2005

L4-4



Active Areas of Research

◇ Functional Genomics:

- Understand how an interactive network of genes affect a chain of metabolic pathways to ultimately determine the phenotypes

◇ Comparative Genomics:

- Relate genes within and across species to understand their evolutionary relationship...Phylogeny.

10/18/2005

© Bud Mishra, 2005

L4-5



Active Areas of Research

◇ Systems Biology:

- Interaction between proteins (membrane and soluble ones) to determine the dynamics of a cell.
- Interaction among a heterogeneous population of cells.

◇ Rational Drug Design:

- Design of drugs and delivery systems to modify the dynamics of the cells.

10/18/2005

© Bud Mishra, 2005

L4-6



Some Biology

10/18/2005

© Bud Mishra, 2005

L4-7



Introduction to Biology

- ◇ **Genome:**
 - Hereditary information of an organism is encoded in its DNA and enclosed in a cell (unless it is a virus). All the information contained in the DNA of a single organism is its genome.
- ◇ **DNA molecule**
 - can be thought of as a very long sequence of nucleotides or bases:
 - ❖ $\Sigma = \{A, T, C, G\}$

10/18/2005

© Bud Mishra, 2005

L4-8



Complementarity

- ◇ DNA is a double-stranded polymer
 - should be thought of as a pair of sequences over Σ .
- ◇ A relation of complementarity
 - A, T, C, G
 - If there is an A (resp., T, C, G) on one sequence at a particular position then the other sequence must have a T (resp., A, G, C) at the same position.
- ◇ The sequence length
 - Is measured in terms of base pairs (bp): Human (*H. sapiens*) DNA is 3.3×10^9 bp, about 6 ft of DNA polymer completely stretched out!

10/18/2005

© Bud Mishra, 2005

L4-9



Genome Size

- ◇ The genomes vary widely in size:
 - Few thousand base pairs for viruses to 2×10^{11} bp for certain amphibian and flowering plants.
 - Coliphage MS2 (a virus) has the smallest genome: only 3.5×10^3 bp.
 - Mycoplasmas (a unicellular organism) has the smallest cellular genome: 5×10^5 bp.
 - *C. elegans* (nematode worm, a primitive multicellular organism) has a genome of size $\approx 10^8$ bp.

Species	Haploid Genome Size	Chromosome Number
<i>E. coli</i>	4.64×10^6	1
<i>S. cerevisiae</i>	1.205×10^7	16
<i>C. elegans</i>	10^8	11/12
<i>D. melanogaster</i>	1.7×10^8	4
<i>M. musculus</i>	3×10^9	20
<i>H. sapiens</i>	3×10^9	23
<i>A. Cepa (Onion)</i>	1.5×10^{10}	8

10/18/2005

© Bud Mishra, 2005

L4-10



DNA) Structure and Components

- ◇ **Double helix**
 - The usual configuration of DNA is in terms of a **double helix** consisting of two **chains** or **strands** coiling around each other with two alternating grooves of slightly different spacing.
 - The "backbone" in each strand is made of alternating sugar molecules (Deoxyribose residues: $C_5 O_4 H_{10}$) and phosphate ($(P O_4)^{-3}$) molecules.
- ◇ **Each of the four bases, an almost planar nitrogenic organic compound, is connected to the sugar molecule.**
 - The bases are: Adenine) A; Thymine) T; Cytosine) C; Guanine) G

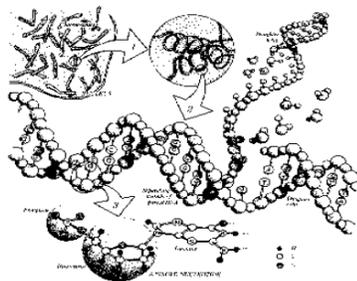
10/18/2005

© Bud Mishra, 2005

L4-11



Genome in Detail



The Human Genome at Four Levels of Detail.

Apart from reproductive cells (gametes) and mature red blood cells, every cell in the human body contains 23 pairs of chromosomes, each a packet of compressed and entwined DNA (1, 2).

10/18/2005

© Bud Mishra, 2005

L4-12



DNA) Structure and Components

- ◇ **Complementary base pairs**
 - (A-T and C-G) are connected by hydrogen bonds and the base-pair forms a coplanar "rung"
 - ❖ Cytosine and thymine are smaller (lighter) molecules, called pyrimidines
 - ❖ Guanine and adenine are bigger (bulkier) molecules, called purines.
 - ❖ Adenine and thymine allow only for double hydrogen bonding, while cytosine and guanine allow for triple hydrogen bonding.

10/18/2005

© Bud Mishra, 2005

L4-13



DNA) Structure and Components

- ◇ **Chemically inert and mechanically rigid and stable**
 - Thus the chemical (through hydrogen bonding) and the mechanical (purine to pyrimidine) constraints on the pairing lead to the complementarity and makes the double stranded DNA both chemically inert and mechanically quite rigid and stable.
- ◇ **Most uninteresting molecule:**
 - "DNA, on its own, does nothing," smirked Natalie Angier recently. "It can't divide, it can't keep itself clean or sit up properly — proteins that surround it do all those tasks. Stripped of context within the body's cells ... DNA is helpless, speechless — DOA."

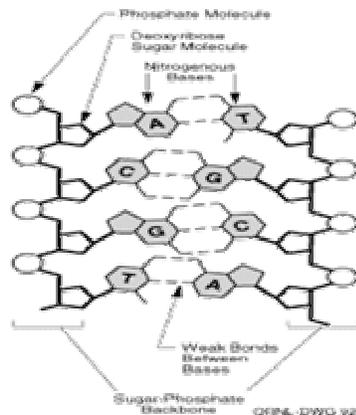
10/18/2005

© Bud Mishra, 2005

L4-14



DNA Structure.



- The four nitrogenous bases of DNA are arranged along the sugar- phosphate backbone in a particular order (the DNA sequence), encoding all genetic instructions for an organism. Adenine (A) pairs with thymine (T), while cytosine (C) pairs with guanine (G). The two DNA strands are held together by weak bonds between the bases.

10/18/2005

© Bud Mishra, 2005

L4-15



DNA) Structure and Components

- ◇ The building blocks of the DNA molecule are four kinds of deoxyribonucleotides,
 - where each deoxyribonucleotide is made up of a sugar residue, a phosphate group and a base.
 - From these building blocks (or related, dNTPs deoxyribonucleoside triphosphates) one can synthesize a strand of DNA.

10/18/2005

© Bud Mishra, 2005

L4-16



DNA) Structure and Components

- ◇ **The sugar molecule**
 - in the strand is in the shape of a pentagon (4 carbons and 1 oxygen) in a plane parallel to the helix axis and with the 5th carbon (5' C) sticking out.
- ◇ **The phosphodiester bond (-O-P-O-)**
 - between the sugars connects this 5' C to a carbon in the pentagon (3' C) and provides a directionality to each strand.
- ◇ **The strands in a double-stranded DNA molecule are antiparallel.**

10/18/2005

© Bud Mishra, 2005

L4-17



The Central Dogma

- ◇ **The central dogma (due to Francis Crick in 1958) states that these information flows are all unidirectional:**
 - "The central dogma states that once 'information' has passed into protein it cannot get out again. The transfer of information from nucleic acid to nucleic acid, or from nucleic acid to protein, may be possible, but transfer from protein to protein, or from protein to nucleic acid is impossible. Information means here the precise determination of sequence, either of bases in the nucleic acid or of amino acid residues in the protein."

10/18/2005

© Bud Mishra, 2005

L4-18



RNA and Transcription

- ◇ **The polymer RNA (ribonucleic acid)**
 - is similar to DNA but differ in several ways:
 - ❖ it's single stranded;
 - ❖ its nucleotide has a ribose sugar (instead of deoxyribose) and
 - ❖ it has the pyrimidine base uracil, U, substituting thymine, T—U is complementary to A like thymine.

10/18/2005

© Bud Mishra, 2005

L4-19



RNA and Transcription

- ◇ **RNA molecule tends to fold back on itself to make helical twisted and rigid segments.**
 - For instance, if a segment of an RNA is
➤ **5' - GGGGAAAACCCC - 3'**,
then the C's fold back on the G's to make a hairpin structure (with a 4bp stem and a 5bp loop).
 - The secondary RNA structure can even be more complicated, for instance, in case of E. coli, Ala tRNA (transfer RNA) forms a cloverleaf shape.
 - Prediction of RNA structure is an interesting computational problem.

10/18/2005

© Bud Mishra, 2005

L4-20



RNA, Genes and Promoters

- ◇ A specific region of DNA that determines the synthesis of proteins (through the transcription and translation) is called a gene
 - Originally, a gene meant something more abstract—a unit of hereditary inheritance.
 - Now a gene has been given a physical molecular existence.
- ◇ Transcription of a gene to a messenger RNA, mRNA,
 - is keyed by an RNA polymerase enzyme, which attaches to a core promoter (a specific sequence adjacent to the gene).

10/18/2005

© Bud Mishra, 2005

L4-21



RNA, Genes and Promoters

- ◇ Regulatory sequences such as silencers and enhancers control the rate of transcription
 - by their influence on the RNA polymerase through a feedback control loop involving many large families of activator and repressor proteins that bind with DNA and
 - which in turn, transpond the RNA polymerase by coactivator proteins and basal factors.

10/18/2005

© Bud Mishra, 2005

L4-22



Transcriptional Regulation

- ◇ The entire structure of transcriptional regulation of gene expression is rather dispersed and fairly complicated:
 - The enhancer and silencer sequences occur over a wide region spanning many Kb's from the core promoter on either directions;
 - A gene may have many silencers and enhancers and can be shared among the genes;

10/18/2005

© Bud Mishra, 2005

L4-23



Transcriptional Regulation

- ◇ The enhancer and silencer sequences
 - They are not unique—different genes may have different combinations;
 - The proteins involved in control of the RNA polymerase number around 50 and
 - Different cliques of transcriptional factors operate in different cliques.
- ◇ Any disorder in their proper operation can lead to cancer, immune disorder, heart disease, etc

10/18/2005

© Bud Mishra, 2005

L4-24



Transcription

- ◇ **The transcription of DNA in to Mrna**
 - is performed with a single strand of DNA (the sense strand) around a gene.
 - This newly synthesized mRNA are capped by attaching special nucleotide sequences to the 5' and 3' ends.
- ◇ **This molecule is called a pre-mRNA.**

10/18/2005

© Bud Mishra, 2005

L4-25



Transcription

- ◇ **The double helix**
 - Untwists momentarily to create a transcriptional bubble which moves along the DNA in the 3' - 5' direction (of the sense strand)
 - As the complementary mRNA synthesis progresses adding one RNA nucleotide at a time at the 3' end of the RNA, attaching an U (respectively, A, G and C) for the corresponding DNA base of A (respectively, T, C and G),
 - Ending when a termination signal (a special sequence) is encountered.

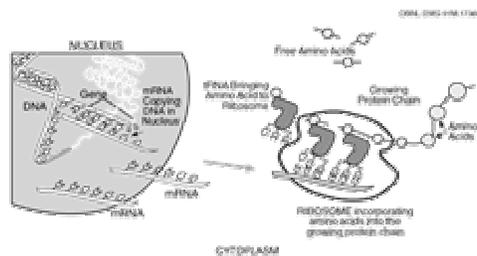
10/18/2005

© Bud Mishra, 2005

L4-26



Gene Expression



◊When genes are expressed, the genetic information (base sequence) on DNA is first **transcribed** (copied) to a molecule of messenger RNA, **mRNA**.

◊The mRNAs leave the cell nucleus and enter the cytoplasm, where triplets of bases (**codons**) forming the genetic code specify the particular amino acids that make up an individual protein.

◊This process, called **translation**, is accomplished by **ribosomes** (cellular components composed of proteins and another class of RNA) that read the genetic code from the mRNA, and transfer RNAs (**tRNAs**) that transport amino acids to the ribosomes for attachment to the growing protein.

10/18/2005

© Bud Mishra, 2005

L4-27



Interrupted Genes

◊ Exons and Introns

- In eukaryotic cells, the region of DNA transcribed into a pre-mRNA involves more than just the information needed to synthesize the proteins.
- The DNA containing the code for protein are the exons, which are interrupted by the introns, the non-coding regions.

10/18/2005

© Bud Mishra, 2005

L4-28



Exons and Introns

- ◇ **Thus pre-mRNA**
 - contains both exons and introns and is altered to excise all the intronic subsequences in preparation for the translation process—this is done by the spliceosome.
- ◇ **The location of splice sites,**
 - separating the introns and exons, is dictated by short sequences and simple rules such as
 - “introns begin with the dinucleotide GT and end with the dinucleotide AG” (the GT-AG rule).

10/18/2005

© Bud Mishra, 2005

L4-29



Protein and Translation

- ◇ **The translation process**
 - begins at a particular location of the mRNA called the translation start sequence (usually AUG) and is mediated by the transfer RNA (tRNA), made up of a group of small RNA molecules, each with specificity for a particular amino acid.
- ◇ **The tRNA's**
 - carry the amino acids to the ribosomes, the site of protein synthesis, where they are attached to a growing polypeptide.

10/18/2005

© Bud Mishra, 2005

L4-30



Protein and Translation

- ◇ **The translation stops**
 - when one of the three trinucleotides UAA, UAG or UGA is encountered.
- ◇ **Codon**
 - Each 3 consecutive (nonoverlapping) bases of mRNA (corresponding to a codon) codes for a specific amino acid.
 - There are $4^3 = 64$ possible trinucleotide codons belonging to the set
 - » $\{U, A, G, C\}^3$

10/18/2005

© Bud Mishra, 2005

L4-31



Genetic Codes

- ◇ **Redundancy in Codons:**
 - The codon AUG is the start codon and the codons UAA, UAG and UGA are the stop codons.
 - That leaves 60 codons to code for 20 amino acids with an expected redundancy of 3!
 - Multiple codons (one to six) are used to code a single amino acid.
- ◇ **Open reading frame (ORF)**
 - The line of nucleotides between and including the start and stop codons.

10/18/2005

© Bud Mishra, 2005

L4-32



ORF

- ◇ All the information of interest to us resides in the ORF's.
- ◇ The mapping from the codons to amino acid (and naturally extended to a mapping from ORF's polypeptides by a homomorphism) given by

$$F_p : \{U, A, G, C\}^3$$

$$\rightarrow \{A, R, D, N, C, E, Q, G, H, I, L, K, M, F, P, S, T, W, Y, V\}$$

10/18/2005

© Bud Mishra, 2005

L4-33



Amino Acids with Codes

A	Ala	alanine	GC(U+A+C+G)
C	Cys	cysteine	UG(U+C)
D	Asp	aspartic acid	GA(U+C)
E	Glu	glutamic acid	GA(G+A)
F	Phe	phenylalanine	UU(U+C)
G	Gly	glycine	GG(U+A+C+G)
H	His	histine	CA(U+C)
I	Ile	Isoleucine	AU(U+A+C)
K	Lys	lysine	AA(A+G)
L	Leu	leucine	(C+U)U(A+G) + CU(U+C)
M	Met	methionine	AUG
N	Asn	asparagine	AA(U+C)
P	Pro	proline	CC(U+A+C+G)
Q	Gln	glutamine	CA(A+G)
R	Arg	arginine	(A+C)G(A+G)+CG(U+C)
S	Ser	serine	(AG+UC)(U+C)+UC(A+G)
T	Thr	threonine	AC(U+A+C+G)
V	Val	valine	GU(U+A+C+G)
W	Trp	tryptophan	UGG
Y	Tyr	tyrosine	UA(U+C)

10/18/2005

© Bud Mishra, 2005

L4-34



The Cell

- ◇ **Small coalition of a set of genes**
 - held together in a set of chromosomes (and even perhaps unrelated extrachromosomal elements).
- ◇ **Set of machinery**
 - made of proteins, enzymes, lipids and organelles taking part in a dynamic process of information processing.

10/18/2005

© Bud Mishra, 2005

L4-35



The Cell

- ◇ **In eukaryotic cells**
 - the genetic materials are enclosed in the cell nucleus separated from the other organelles in the cytoplasm by a membrane.
- ◇ **In prokaryotic cells**
 - the genetic materials are distributed homogeneously as it does not have a nucleus.
 - Example of prokaryotic cells are bacteria with a considerably simple genome.

10/18/2005

© Bud Mishra, 2005

L4-36



Organelles

- ◇ **The organelles common to eukaryotic plant and animal cells include**
 - **Mitochondria** in animal cells and chloroplasts in plant cells (responsible for energy production);
 - A **Golgi apparatus** (responsible for modifying, sorting and packaging various macromolecules for distribution within and outside the cell);
 - **Endoplasmic reticulum** (responsible for synthesizing protein); and
 - **Nucleus** (responsible for holding the DNA as chromosomes and replication and transcription).

10/18/2005

© Bud Mishra, 2005

L4-37



Chromosomes

- ◇ **The entire cell**
 - is contained in a sack made of plasma membrane.
 - In plant cells, they are further surrounded by a cellulose cell wall.
- ◇ **The nucleus of the eukaryotic cells**
 - contain its genome in several chromosomes, where each chromosome is simply a single molecule of DNA as well as some proteins (primarily histones).

10/18/2005

© Bud Mishra, 2005

L4-38



Chromosomes

- ◇ **The chromosomes**
 - can be a **circular** or **linear**, in which case the ends are capped with special sequence of telomeres.
- ◇ **The protein**
 - in the nucleus binds to the DNA and effects the compaction of the very long DNA molecules.
- ◇ **Ploidy**
 - In somatic cells of most eukaryotic organisms, the chromosomes occur in homologous pairs,
 - Exceptions: X and Y –sex chromosomes.

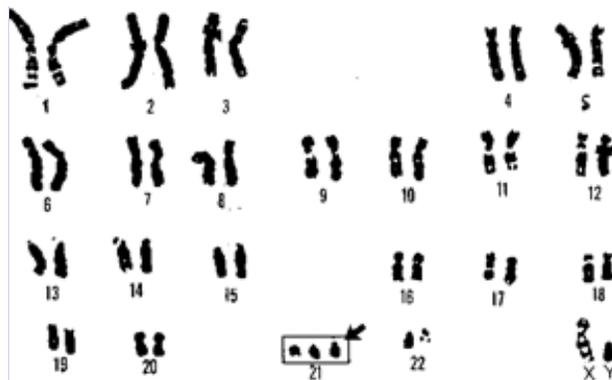
10/18/2005

© Bud Mishra, 2005

L4-39



Chromosomes



◇Karyotype.

◇Microscopic examination of chromosome size and banding patterns identifies 24 different chromosomes in a karyotype, which is used for diagnosis of genetic diseases.

◇The extra copy of chromosome 21 (trisomy) in this karyotype implies Down's syndrome.

10/18/2005

© Bud Mishra, 2005

L4-40



Ploidy

- ◇ **Gametes contain only unpaired chromosomes;**
 - the egg cell contains only X chromosome and the sperm cell either an X or an Y chromosome. The male has X and Y chromosomes; the female, 2 X's.
 - Cells with single unpaired chromosomes are called haploid;
 - Cells with homologous pairs, diploid;
 - Cells with homologous triplet, quadruplet, etc., chromosomes are called polyploid—many plant cells are polyploid.

10/18/2005

© Bud Mishra, 2005

L4-41



Chromosomal Aberrations

- ◇ **Point mutations**
- ◇ **Breakage**
- ◇ **Translocation (Among non-homologous chromosomes.)**
- ◇ **Formation of acentric and dicentric chromosomes.**
- ◇ **Gene Conversions**
- ◇ **Amplification and deletions**
- ◇ **Jumping genes a Transposition of DNA segments**
- ◇ **Programmed rearrangements a E.g., antibody responses.**

10/18/2005

© Bud Mishra, 2005

L4-42



Point Mutations

◇ In exon:

- Can change the protein,
 - ❖ if it is transcriptional factor, it can affect many other genes
- Can terminate the mRNA too early by changing a non-stop codon into a stop codon
 - (NMD: Nonsense Mediated Degradation)
- Small indel can cause frame-shift, thus changing the entire protein

10/18/2005

© Bud Mishra, 2005

L4-43



Point Mutations

◇ In promoter:

- Can change its regulation: Over express, under express or silence

◇ In intron:

- Can change the splicing patterns

10/18/2005

© Bud Mishra, 2005

L4-44



Loss or gain

- ◇ **Translocation:**
 - Can fuse two genes
 - Can activate a silent gene by placing it near some active regulatory region
- ◇ **Amplification:**
 - Over expression; Highly active gene
- ◇ **Deletion**
 - Under expression: Inactive gene

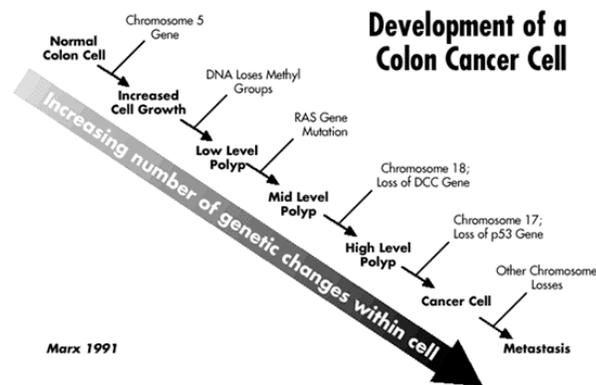
10/18/2005

© Bud Mishra, 2005

L4-45



Amplifications & Deletions



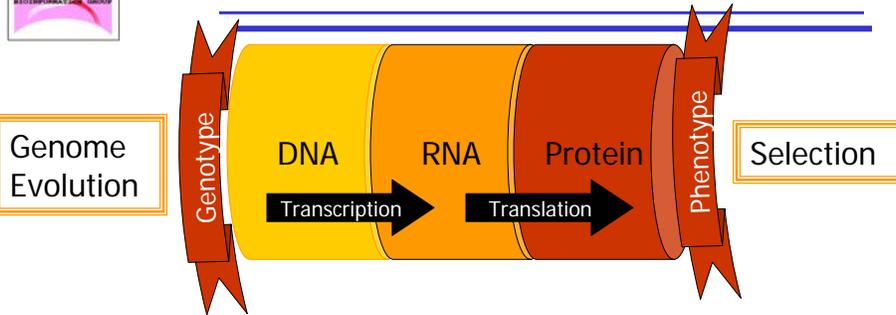
10/18/2005

© Bud Mishra, 2005

L4-46



The New Synthesis



Part-lists, Annotation, Ontologies

10/18/2005

© Bud Mishra, 2005

L4-47



Cancer Initiation and Progression

Mutations, Translocations, Amplifications, Deletions
Epigenomics (Hyper & Hypo-Methylation)
Alternate Splicing

Cancer Initiation and Progression

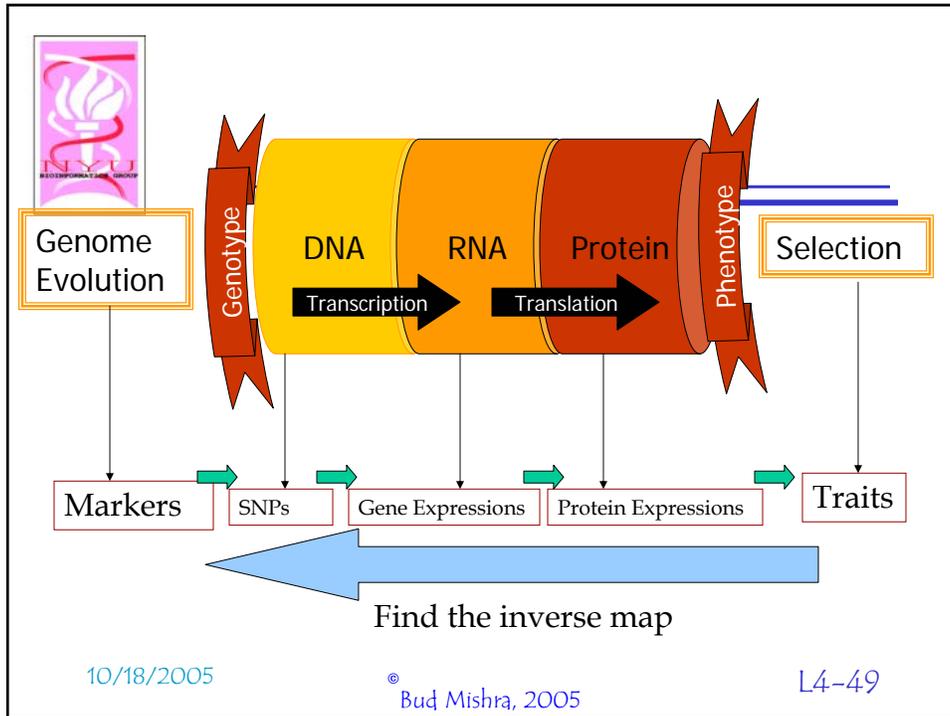


Proliferation, Motility, Immortality, Metastasis, Signaling

10/18/2005

© Bud Mishra, 2005

L4-48



 To be continued...

 ...

10/18/2005 © Bud Mishra, 2005 L4-50